# Open-Source Projects: Machine Learning for Transportation Data Imputation and Prediction

Reproducible Research Workshop

TRB 103rd Annual Meeting · Washington, D.C., USA

**Nicolas Saunier & Xinyu Chen**

Polytechnique Montreal, Canada

January 11, 2024

Open-source & reproducible research:

❶ GitHub: https://github.com/xinychen

❷ Slides: https://xinychen.github.io/slides/transdim.pdf

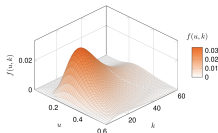❸ Project website: https://spatiotemporal-data.github.io

ML algorithms



transdim

(1.1k stars)

Visualization tools



awesome-latex-drawing

(1.2k stars)

Other projects at Polytechnique Montreal in video processing for user behavior and safety analysis https://trafficintelligence.confins.net

# Why?

Academia:

- Open research environment (w.r.t. our team & followers)
- Push by funding agencies and academic institutions: "Research Data Management Policy"
- Interact with researchers from different fields
- Provide platform and benchmark for comparison
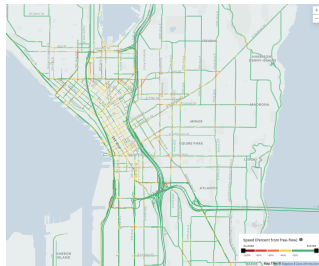- Stimulate new algorithmic ideas

Industry:

- Provide solutions to some realistic data problems
- Easy to implement and produce results

# Storytelling with Data

- Uber (hourly) movement speed data
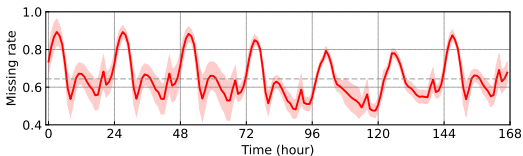


NYC movement



Seattle movement

- {road segment, time step (hour), average speed}
- $\boldsymbol{Y} \in \mathbb{R}^{N \times T}$ with $N$ spatial locations $\times$ $T$ time steps
- Computing hourly speed: Road segments have $5+$ unique trips.

**Issue**: Insufficient sampling of ridesharing vehicles on the road network!
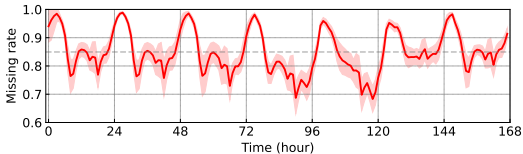
# Storytelling with Data

**High-dimensional & sparse**

- **NYC** movement speed data (2019)
  - 98,210 road segments & 8,760 time steps (hours)
  - Overall missing rate: 64.43%



- **Seattle** movement speed data (2019)
  - 63,490 road segments & 8,760 time steps (hours)
  - Overall missing rate: 84.95%



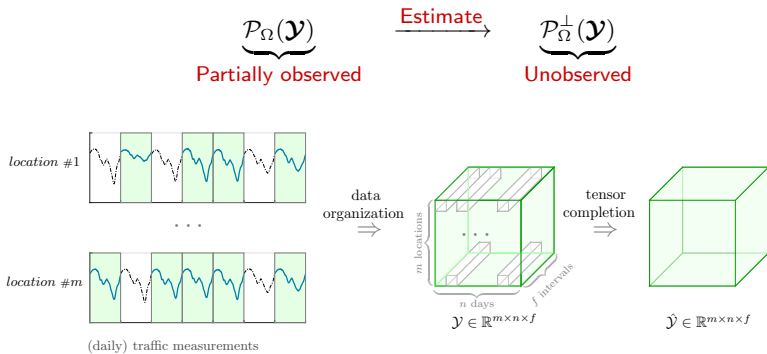- How to quantify data quality? Address sparsity?

# Reformulate Traffic Data Imputation

**Estimation:**

- Imputation & interpolation & forecasting
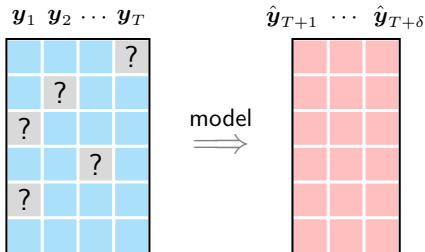
Imputing missing traffic data

- Tensor completion (Observed index set $\Omega$)

$$\underbrace{\mathcal{P}_\Omega(\boldsymbol{\mathcal{Y}})}_{\text{Partially observed}} \xrightarrow{\text{Estimate}} \underbrace{\mathcal{P}_\Omega^\perp(\boldsymbol{\mathcal{Y}})}_{\text{Unobserved}}$$



$location\#1$

$\cdots$

$location\#m$

(daily) traffic measurements

data organization $\Rightarrow$

$m$ locations

$\cdots$

$n$ days

$f$ intervals

$\mathcal{Y} \in \mathbb{R}^{m \times n \times f}$

tensor completion $\Rightarrow$

$\hat{\mathcal{Y}} \in \mathbb{R}^{m \times n \times f}$

# Reformulate Traffic Forecasting

Forecasting urban traffic states with sparse data

- Problem definition (δ-step ahead forecasting)

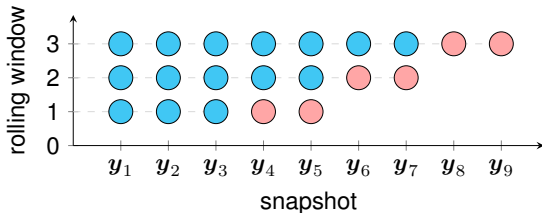$$\underbrace{\{\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_T\}}_{\text{Current traffic states}} \xrightarrow{\text{Estimate}} \underbrace{\{\hat{\boldsymbol{y}}_{T+1}, \hat{\boldsymbol{y}}_{T+2}, \ldots, \hat{\boldsymbol{y}}_{T+\delta}\}}_{\text{Future traffic states}}$$

# Reformulate Traffic Forecasting

(Rolling) Forecasting urban traffic states with sparse data

1st rolling step: $\{\boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{y}_3\} \to \{\boldsymbol{y}_4, \boldsymbol{y}_5\}$

2nd rolling step: $\{\boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{y}_3, \boldsymbol{y}_4, \boldsymbol{y}_5\} \to \{\boldsymbol{y}_6, \boldsymbol{y}_7\}$

3rd rolling step: $\underbrace{\{\boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{y}_3, \boldsymbol{y}_4, \boldsymbol{y}_5, \boldsymbol{y}_6, \boldsymbol{y}_7\}}_{\text{Current traffic states}} \to \underbrace{\{\boldsymbol{y}_8, \boldsymbol{y}_9\}}_{\text{Future traffic states}}$

# Python Implementation

**Tools & Packages**

| Python | CPU computing | GPU computing |
|--------|---------------|---------------|



*NumPy for GPU

Traffic Data Processing

- Data format: `.npz` (compressed format)

- Easy to use
  - Connect with `numpy` (for **CPU**)
  - Connect with `cupy` (for **GPU**)

NYC Uber movement dataset[1]:

- e.g., `hourly_speed_mat_2019_1.npz` (91 MB)
  - $98210 \times 744$ matrix
  - 23,228,581 observations
  - Whole month of January 2019

---

[1]https://github.com/xinychen/tracebase

# Python Implementation

## Pre-process several open datasets[2]

### Open data

In this project, we have adapted some publicly available data sets into our experiments. The original links for these data are summarized as follows,

- **Multivariate time series**
  - Birmingham parking data set
  - California PeMS traffic speed data set (large-scale)
  - Guangzhou urban traffic speed data set
  - Hangzhou metro passenger flow data set
  - London urban movement speed data set (other cities are also available at Uber movement project)
  - Portland highway traffic data set (including traffic volume/speed/occupancy, see data documentation)
  - Seattle freeway traffic speed data set
- **Multidimensional time series**
  - New York City (NYC) taxi data set
  - Pacific surface temperature data set

For example, if you want to view or use these data sets, please download them at the ../datasets/ folder in advance, and then run the following codes in your Python console:

```
import scipy.io

tensor = scipy.io.loadmat('../datasets/Guangzhou-data-set/tensor.mat')
tensor = tensor['tensor']
```

[2] https://github.com/xinychen/transdim

# Python Implementation

Machine learning algorithms:

- Low-rank tensor completion
- Temporal matrix factorization

```python
def LRTC(dense_tensor, sparse_tensor, alpha, rho, theta, epsilon, maxiter):
    """Low-Rank Tensor Completion with Truncated Nuclear Norm, LRTC-TNN."""

    dim = np.array(sparse_tensor.shape)
    pos_missing = np.where(sparse_tensor == 0)
    pos_test = np.where((dense_tensor != 0) & (sparse_tensor == 0))
    dense_test = dense_tensor[pos_test]
    del dense_tensor

    X = np.zeros(np.insert(dim, 0, len(dim))) # \boldsymbol{\mathcal{X}}
    T = np.zeros(np.insert(dim, 0, len(dim))) # \boldsymbol{\mathcal{T}}
    Z = sparse_tensor.copy()
    last_tensor = sparse_tensor.copy()
    snorm = np.sqrt(np.sum(sparse_tensor ** 2))
    it = 0
    while True:
        rho = min(rho * 1.05, 1e5)
        for k in range(len(dim)):
            X[k] = mat2ten(svt_tnn(ten2mat(Z - T[k] / rho, k), alpha[k] / rho, np.int(np.ceil(
            Z[pos_missing] = np.mean(X + T / rho, axis = 0)[pos_missing]
            T = T + rho * (X - np.broadcast_to(Z, np.insert(dim, 0, len(dim))))
            tensor_hat = np.einsum('k, kmnt -> mnt', alpha, X)
            tol = np.sqrt(np.sum((tensor_hat - last_tensor) ** 2)) / snorm
            last_tensor = tensor_hat.copy()
            it += 1
            if (it + 1) % 50 == 0:
                print('Iter: {}'.format(it + 1))
                print('MAPE: {:.6}'.format(compute_mape(dense_test, tensor_hat[pos_test])))
                print('RMSE: {:.6}'.format(compute_rmse(dense_test, tensor_hat[pos_test])))
                print()
            if (tol < epsilon) or (it >= maxiter):
                break

    print('Imputation MAPE: {:.6}'.format(compute_mape(dense_test, tensor_hat[pos_test])))
    print('Imputation RMSE: {:.6}'.format(compute_rmse(dense_test, tensor_hat[pos_test])))
    print()

    return tensor_hat
```

- Define nonstationary temporal matrix factorization ( `notmf` ).

```python
def notmf(dense_mat, sparse_mat, rank, d, lmbda, rho, season, maxiter):
    dim1, dim2 = sparse_mat.shape
    W = 0.01 * np.random.randn(rank, dim1)
    X = 0.01 * np.random.randn(rank, dim2)
    A = 0.01 * np.random.randn(rank, d * rank)
    if np.isnan(sparse_mat).any() == False:
        ind = sparse_mat != 0
        pos_test = np.where((dense_mat != 0) & (sparse_mat == 0))
    elif np.isnan(sparse_mat).any() == True:
        pos_test = np.where((dense_mat != 0) & (np.isnan(sparse_mat)))
        ind = ~np.isnan(sparse_mat)
        sparse_mat[np.isnan(sparse_mat)] = 0
    dense_test = dense_mat[pos_test]
    del dense_mat
    Psi = generate_Psi(dim2, d, season)
    show_iter = 100
    temp = np.zeros((d * rank, dim2 - d - season))
    for it in range(maxiter):
        W = conj_grad_w(sparse_mat, ind, W, X, rho)
        X = conj_grad_x(sparse_mat, ind, W, X, A, Psi, d, lmbda, rho)
        for k in range(d, d + 1):
            temp[(k - 1) * rank : k * rank, :] = X @ Psi[k].T
        A = X @ Psi[0].T @ np.linalg.pinv(temp)
        mat_hat = W.T @ X
        if (it + 1) % show_iter == 0:
            temp_hat = mat_hat[pos_test]
            print('Iter: {}'.format(it + 1))
            print('MAPE: {:.6}'.format(compute_mape(dense_test, temp_hat)))
            print('RMSE: {:.6}'.format(compute_rmse(dense_test, temp_hat)))
            print()
    return mat_hat, W, X, A
```

# Python Implementation

Easy to **switch from CPU to GPU**



```
import numpy as np
```

$\Rightarrow$



```
import cupy as np
```

Only use `numpy`? Advantage:

- Fewer packages can improve the reproducibility

# Post Something That Matters

Post well-documented data processing files (e.g., processing Chicago taxi data)

- Beginners to build coding skills
- Researchers to build research ideas



Source: https://spatiotemporal-data.github.io/Chicago-mobility/taxi-data

# Post Something That Matters

Post <span style="color:red">scientific problems</span> (e.g., spatiotemporal data modeling)

### Optimizing Interpretable Time-Varying Autoregression with Orthogonal Constraints

Generally speaking, any spatiotemporal data in the form of a matrix can be written as $\boldsymbol{Y} \in \mathbb{R}^{N \times T}$ with $N$ spatial areas/locations and $T$ time steps. To discover interpretable spatial/temporal patterns, one can build a time-varying autoregression on the time snapshots $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_T \in \mathbb{R}^N$ (Chen et al., 2023). The time-varying coefficients in the autoregression allow one to characterize the time-varying system behavior, but the challenges still remain.

To capture interpretable modes/patterns, one can use tensor factorization formulas to parameterize the coefficients and the optimization problem can be easily built. However, a great challenge would be how to make the modes "more interpretable", specifically, e.g., how to learn orthogonal modes in the modeling process. In this post, we present an optimization problem of the time-varying autoregression with orthogonal constraints as follows,

$$\min_{\boldsymbol{W}, \boldsymbol{G}, \boldsymbol{V}, \boldsymbol{X}} \frac{1}{2} \sum_{t=2}^{T} \left\| \boldsymbol{y}_t - \boldsymbol{W} \boldsymbol{G} (\boldsymbol{x}_t^\top \otimes \boldsymbol{V})^\top \boldsymbol{y}_{t-1} \right\|_2^2$$
$$\text{s.t.} \begin{cases} \boldsymbol{W}^\top \boldsymbol{W} = \boldsymbol{I}_R \\ \boldsymbol{V}^\top \boldsymbol{V} = \boldsymbol{I}_R \\ \boldsymbol{X}^\top \boldsymbol{X} = \boldsymbol{I}_R \end{cases}$$

where $\boldsymbol{W} \in \mathbb{R}^{N \times R}$ and $\boldsymbol{X} \in \mathbb{R}^{(T-1) \times R}$ refer to as the spatial modes and the temporal modes, respectively. This model can discover urban mobility transition patterns.

Source: <span style="color:blue">https://spatiotemporal-data.github.io/probs/orth-var</span>

# Reproducibility Challenges

Resources

- Time / man power to create and maintain the material
- Storage (for data)
- Sustainable platforms

Creating and fostering a community

- Example of Traffic Intelligence, open source software for traffic video processing and safety analysis

The devil is in the details

- Have you tried to compare a new video object detector to the state of the art?

# Thanks for your attention!

## Any Questions?

**About Nicolas Saunier**:
- 🏠 Homepage: https://nicolas.saunier.confins.net
- 🐙 GitHub: https://github.com/nsaunier
- ✉ How to reach me: nicolas.saunier@polymtl.ca

**About Xinyu Chen**:
- 🏠 Homepage: https://xinychen.github.io
- 🐙 GitHub: https://github.com/xinychen
- ✉ How to reach me: chenxy346@gmail.com